

# Stochastic Backpropagation

Presented By Sahil Manocha

Danilo J. Rezende, Shakir Mohamed, Daan Wierstra

2014

# Exponential Families

$$p(x|\eta) = h(x)g(\eta)\exp\{\eta^T u(x)\}$$

## Maximum Likelihood Estimation

$$p(X|\eta) = \int h(x)g(\eta)\exp\{\eta^T u(x)\}$$

Differentiating, we get

$$\begin{aligned}\frac{\partial p(X|\eta)}{\partial \eta} &= \nabla g(\eta) \int h(x) \exp\{\eta^T u(x)\} dx \\ &\quad + g(\eta) \int h(x) \exp\{\eta^T u(x)\} u(x) dx\end{aligned}$$

# Exponential Families

Alternatively,  $\eta = \eta(\theta)$ ,

$$p(x|\eta) = h(x)\exp\{\eta^T u(x) - A(\eta)\}$$

$A$  is the *log-partition function*.

# Exponential Families

## Theorem

The log-partition function  $\theta \rightarrow A(\theta)$  is infinitely differentiable on its open domain  $D := \{\theta \in \mathbb{R}^d : A(\theta) < \infty\}$ . Moreover,  $A$  is convex.

## Proof.

For convexity, let  $\theta_\lambda = \lambda\theta_1 + (1 - \lambda)\theta_2$ , where  $\theta_1, \theta_2 \in D$ . Then,  $\frac{1}{\lambda} \geq 1$  and  $\frac{1}{1-\lambda} \geq 1$ , and Holder's inequality is applicable. (Since the coefficients are conjugate exponents). We get,

$$\begin{aligned} & \log \int h(x) \exp(\langle \theta_\lambda, u(x) \rangle) dx \\ &= \log \int h(x) \exp(\langle \theta_1, u(x) \rangle)^\lambda \exp(\langle \theta_2, u(x) \rangle)^{1-\lambda} dx \\ &\leq \log \left( \int h(x) \exp(\langle \theta_1, u(x) \rangle)^\lambda dx \right)^\lambda \left( \int \exp(\langle \theta_2, u(x) \rangle)^{\frac{1-\lambda}{1-\lambda}} dx \right)^{1-\lambda} \end{aligned}$$

# Exponential Distribution

Proof (Cont.)

$$= \lambda \log \int h(x) \exp(\langle \theta_1, u(x) \rangle) dx + (1 - \lambda) \log \int h(x) \exp(\langle \theta_2, u(x) \rangle) dx$$



# Exponential Distribution

Convexity makes estimation in exponential families substantially easier. Indeed, given a sample  $X_1, \dots, X_n$  assume that we estimate  $\theta$  by maximizing likelihood (equivalently, minimizing the log loss):

$$\min_{\theta} \sum_{i=1}^n \log \frac{1}{p_{\theta}(X_i)} = \sum_{i=1}^n [-\langle \theta, u(X_i) \rangle + A(\theta)]$$

which is convex in  $\theta$ .

# Stochastic Backpropagation

Gradient descent methods in latent variable models require computations

$$\nabla_{\theta} \mathbb{E}_{q_{\theta}}[f(x)]$$

$$\theta \sim q_{\theta}(\cdot)$$

$f$  = loss function

Quantity is difficult to compute:

1. expectation is unknown
2. indirect dependency on  $q$

# Bonnet's Theorem

Let  $f(x): \mathbb{R}^d \rightarrow \mathbb{R}$  be a integrable and twice differentiable function. The gradient of the expectation of  $f(x)$  under a Gaussian distribution  $\mathcal{N}(x|\mu, \Sigma)$  with respect to the mean  $\mu$  can be expressed as the expectation of the gradient of  $f(x)$ .

$$\nabla_{\mu_i} \mathbb{E}_{\mathcal{N}(\mu, \Sigma)}[f(x)] = \mathbb{E}_{\mathcal{N}(\mu, \Sigma)}[\nabla_{x_i} f(x)]$$



# Price's Theorem

Let  $f(x): \mathbb{R}^d \rightarrow \mathbb{R}$  be a integrable and twice differentiable function. The gradient of the expectation of  $f(x)$  under a Gaussian distribution  $\mathcal{N}(x|\mu, \Sigma)$  with respect to the covariance  $\Sigma$  can be expressed in terms of the expectation of the Hessian of  $f(x)$  as:

$$\nabla_{\Sigma_{i,j}} \mathbb{E}_{\mathcal{N}(\mu, \Sigma)}[f(x)] = \frac{1}{2} \mathbb{E}_{\mathcal{N}(\mu, \Sigma)}[\nabla_{x_i, x_j}^2 f(x)]$$

# Backpropagation

By applying chain rule, we get:

$$\nabla_{\theta} \mathbb{E}_{\mathcal{N}(\mu, \Sigma)}[f(x)] = \mathbb{E}_{\mathcal{N}(\mu, \Sigma)} \left[ \mathbf{g}^T \frac{\partial \mu}{\partial \theta} + \frac{1}{2} \text{Tr} \left( H \frac{\partial \Sigma}{\partial \theta} \right) \right]$$

$\mathbf{g}$  : gradient

$H$ : hessian

# For General Distributions

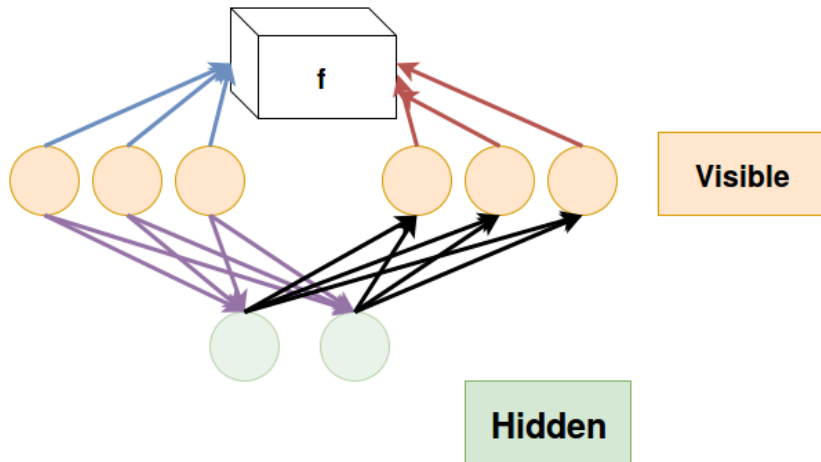
$$\nabla_{\theta} \mathbb{E}_p[f(x)] = \mathbb{E}_{p(x|\theta)}[B(x) \nabla_x f(x)]$$

where  $B$  is a non-linear function.

For exponential distributions:

$$B(x) = \frac{[\nabla_{\theta} \eta u(x) - \nabla_{\theta} A]}{[\nabla_x \log[h(x)] + \eta^T \nabla_x u(x)]}$$

## Relation to RBM



## Energy Function

$$E(\mathbf{v}, \mathbf{h}) = \sum_{i=1}^V \frac{(v_i - b_i)^2}{2\sigma_i^2} - \mathbf{c}^T \mathbf{h} - \sum_{j=1}^V \sum_{i=1}^H \frac{v_i}{\sigma_i} h_j w_{ij}$$

## Probability

$$P(\mathbf{v}) = \sum_{\mathbf{h}} \frac{1}{Z} e^{-E(\mathbf{v}, \mathbf{h})} = \sum_{\mathbf{h}} \frac{1}{Z} e^{-\sum_{i=1}^V \frac{(v_i - b_i)^2}{2\sigma_i^2} + \mathbf{c}^T \mathbf{h} + \sum_{j=1}^V \sum_{i=1}^H \frac{v_i}{\sigma_i} h_j w_{ij}}$$

$$P(\mathbf{v}) = \frac{1}{Z} e^{-F(\mathbf{v})}$$

where  $F(\mathbf{v})$  is free energy

# GB-RBM

## Free Energy

$$F(v) = -\log\left(\sum_{\mathbf{h}} e^{-\sum_{i=1}^V \frac{(v_i - b_i)^2}{2\sigma_i^2} + \mathbf{c}_{\mathbf{h}}^T \mathbf{h} + \sum_{j=1}^H \sum_{i=1}^V \frac{v_i}{\sigma_i} h_j w_{ij}}\right)$$

Simplifying the term within the  $\log$

$$\begin{aligned} & \sum_{\mathbf{h}} e^{-\sum_{i=1}^V \frac{(v_i - b_i)^2}{2\sigma_i^2} + \mathbf{c}_{\mathbf{h}}^T \mathbf{h} + \sum_{j=1}^H \sum_{i=1}^V \frac{v_i}{\sigma_i} h_j w_{ij}} \\ &= e^{-\sum_{i=1}^V \frac{(v_i - b_i)^2}{2\sigma_i^2}} \times \prod_j (e^{c_j + \sum_{i=1}^V \frac{v_i}{\sigma_i} w_{ij}} + 1) \end{aligned}$$

Substituting, we get

$$F(v) = \sum_{i=1}^V \frac{(v_i - b_i)^2}{2\sigma_i^2} - \sum_j \log(e^{c_j + \sum_{i=1}^V \frac{v_i}{\sigma_i} w_{ij}} + 1)$$

# Conclusions

Stochastic Backpropagation possible as shown above.

# References I



Danilo Jimenez Rezende, Shakir Mohamed, Daan Wierstra  
Stochastic Backpropagation and Approximate Inference in  
Deep Generative Models